

All That Glitters Is Not Gold: Four Maturity Stages of Process Discovery Algorithms

Jan Martijn E. M. van der Werf Artem Polyvyanyy
Bart R. van Wensveen Matthieu Brinkhuis
Hajo A. Reijers

Abstract

A process discovery algorithm aims to construct a process model that represents the real-world process stored in event data well; it is precise, generalizes the data correctly, and is simple. At the same time, it is reasonable to expect that better-quality input event data should lead to constructed process models of better quality. However, existing process discovery algorithms omit the discussion of this relationship between the inputs and outputs and, as it turns out, often do not guarantee it. We demonstrate the latter claim using several quality measures for event data and discovered process models. Consequently, this paper requests for more rigor in the design of process discovery algorithms, including properties that relate the qualities of the inputs and outputs of these algorithms. We present four incremental maturity stages for process discovery algorithms, along with concrete guidelines for formulating relevant properties and experimental validation. We then use these stages to review several state of the art process discovery algorithms to confirm the need to reflect on how we perform algorithmic process discovery.

1 Introduction

Process mining studies algorithms that extract process-related information from event data, often recorded in event logs as collections of sequences of activities, each encoding a historical process execution [32]. Process discovery is one of the core subareas of process mining. Process discovery studies algorithms that construct models describing the processes that induced the input event logs as closely as possible. One of the challenges of process discovery is that the *true processes* that generated the input event logs are unknown and, thus, must be inferred from their *samples*, that is, event logs [7, 24].

An algorithm is a sequence of computational steps that transform an *input* into an *output* [9]. Different algorithms exhibit different qualities in terms of properties like correctness, finiteness, definiteness, effectiveness, and efficiency.

Such properties allow us to choose an algorithm that fulfills a particular need, such as performing a guaranteed correct computation within the desired time bounds. A process discovery algorithm transforms a given input event log into an output process model. A process discovery algorithm is often finite (terminates after a finite number of computational steps), definite (each computational step is unambiguous), effective (each computational step can be performed correctly in a finite amount of time), and efficient (the fewer or faster computation steps can be executed the better). However, process discovery algorithms treat quality as a *goal* rather than a guarantee. That is, process discovery algorithms are designed to construct a “good” process model from the input event log [32], where the “goodness” of the model is not established by the internals of the algorithm but by external measures, e.g., precision and recall [7, 32, 23].

Recently, we observed that a process discovery algorithm can construct a good process model from an event log and construct an inferior model from an event log that is of better quality than the original log [23]. This observation triggered a desire to review and refine how the quality of a process discovery algorithm is established. In our conference paper [38], we argue that a process discovery algorithm should come with guarantees formulated in terms of the relationships between its inputs and outputs. This article refines the original contributions with a discussion on statistical sampling techniques and their effects. It makes the following contributions:

- It proposes measures for the quality of event logs, both in the presence and absence of a true process. In the former case, we use standard conformance checking measures, while in the latter case we rely on sampling techniques and measures as studied in statistics;
- It discusses requirements to measure the quality of samples with respect to an original event log and shows the effect different sampling techniques have on the proposed quality measures;
- It provides empirical evidence that existing process discovery algorithms can construct good models from event logs and, at the same time, produce poor models from better logs; and
- It proposes four maturity stages for process discovery algorithms that aim to demonstrate the relation between the quality of input event logs and the quality of output process models.

These proposals for assessing the goodness of process discovery algorithms can help to advance the field. Several benchmarks (cf. [2]) have identified process discovery algorithms that “glitter,” that is, algorithms that produce high-quality models on a limited collection of event logs. We argue that such benchmarks should be complemented with formal analysis to provide quality guarantees with the algorithms. We invite the process mining community to contribute to the discussion of the maturity of process discovery algorithms. In addition, we encourage the authors of existing and future process discovery algorithms to establish the proposed guarantees.

The remainder of the paper is structured as follows. The next section argues why process discovery algorithms need to provide guarantees. Then, a statistical approach to establish event log quality is introduced in Section 3. Next, Section 4 presents four stages of maturity for process discovery algorithms, together with empirical evidence that there are algorithms that do not provide such guarantees. Finally, Sections 5 and 6 are devoted to related work, and conclusions, respectively.

2 Setting the Stage

A recent study [23] revealed that for some event logs process discovery algorithms could return inferior models for better quality input event logs. In this section, we reflect on the consequences of this observation and its implications for process discovery algorithms. We start with introducing process discovery in Section 2.1, and discuss the desired relation between log and model quality for process discovery algorithms in Section 2.2.

2.1 Process Discovery

Process mining projects often start by assuming that some underlying process generates an event log that can be observed, recorded, and used for process discovery. We refer to this underlying entity as the *true process*. The true process is, however, often unknown [7]. Hence, it can only be approximated. Therefore, based on the observed log, process discovery algorithms aim to construct a process model that describes the true process as closely as possible. Formally, given a set of activities A , an event log L is defined as a multiset over finite sequences, called *traces*, over A . A discovery algorithm disc can be described as a relation $\text{disc} \subseteq \mathcal{L}(A) \times 2^{\mathcal{M}(A)}$, where $\mathcal{L}(A)$ and $\mathcal{M}(A)$ are the universe of all possible logs and the universe of all models over A , respectively. Note that some discovery algorithms, for instance the ILP-miner [37], are non-deterministic and can yield different results for the same input log.

To measure how well a discovered process model describes the behavior recorded in the event log, different conformance measures have been proposed [36]. *Precision* is a function $\text{prec} : \mathcal{L}(A) \times \mathcal{M}(A) \rightarrow [0, 1]$ that quantifies the fraction of behavior allowed by the model that was actually observed. *Recall* is a function $\text{rec} : \mathcal{L}(A) \times \mathcal{M}(A) \rightarrow [0, 1]$ that quantifies the observed behavior allowed by the model. For both measures, the value of one denotes perfect conformance between the log and model. For example, precision and recall can be grounded in the notion of topological entropy of the processes described in the model and log [23]. As demonstrated in [23, 29], the entropy-based precision and recall measures satisfy all the requirements for conformance measures proposed by the process mining community [30, 36, 23, 29].

Process discovery algorithms are often designed with a specific quality goal in mind. Several algorithms have *rediscoverability* as their goal: if the unknown, true process that generated the event log has specific properties, and

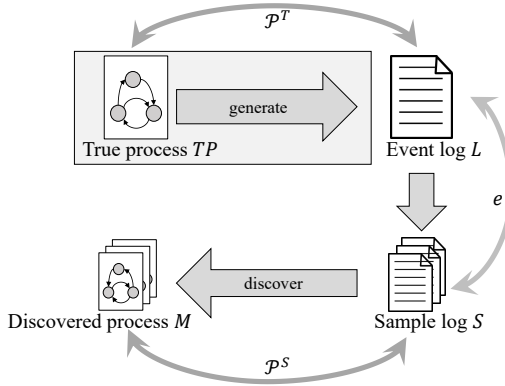


Figure 1: A true process TP generates an event log L with unknown quality \mathcal{P}^T . A sample S drawn from L has some error e . Discovering a model from S results in a process model with quality \mathcal{P}^S .

the event log satisfies certain criteria, then the algorithm ought to discover the true process. For example, the α -miner has the rediscoverability property for structured workflow nets, imposing log completeness as the criterion [35]. Similarly, the Inductive Miner [15] can rediscover process trees under the assumption of activity completeness, i.e., every leaf in the tree should occur at least once in the event log. Another common goal of process discovery algorithms is to construct a model that scores high on one or several conformance measures (e.g., [37, 40, 11]).

2.2 Relating Log Quality and Model Quality

The quality of the results of process discovery algorithms depends on the quality of the input event log. Event logs are often assumed to be faithful representations of the true processes. Let us reflect on the consequences of this assumption. Consider Fig. 1. The true process TP is executed continuously, thus generating a stream of events, from which L is a (non-random) sample [13, 24, 36]. Assume L is a faithful representation of the true process TP . In other words, L has a high model quality \mathcal{P}^T , measured, for example, in terms of precision and recall between L and TP . Therefore, L can be seen as a sample from this stream. Potentially, *samples* of L can be faithful representations of TP as well. Let S be a random sample of L . As it is a random sample, statistical methods can be applied to establish its quality, or lack thereof e , with respect to L . And, because S is an event log itself, it can be used to discover some model M , which has quality \mathcal{P}^S , again measured in terms of conformance measures, but this time between S and M . It follows that if S is a good representation of log L , a process discovery algorithm should construct a model with a quality that approaches \mathcal{P}^T .

Let us draw two samples from L , say S_1 and S_2 . For S_1 , model M_1 is

discovered, with quality \mathcal{P}^{S_1} , and for S_2 a model M_2 is discovered, with quality \mathcal{P}^{S_2} . Suppose S_1 has a higher sample quality than S_2 with respect to L . In other words, S_1 is a better representation of L than S_2 . Intuitively, the quality of M_1 should also be closer to \mathcal{P}^T than the quality of M_2 . In other words, if $e(S_1) \geq e(S_2)$ then one should expect that $\mathcal{P}^{S_1} \geq \mathcal{P}^{S_2}$. Hence, it is desirable that the process discovery algorithm also guarantees that better quality logs result in better quality models.

In real-life situations, the true process that generated the event log is often unknown. In most process mining methods (cf., [6, 39, 31]), the event log is prepared, and then process discovery techniques are applied to unravel a process model. An important concern that these methods do not address relates to the reliability [27] of process mining projects: if the process is repeated on a new observation, i.e., a new event log, to what degree do the results agree between the analyses? Specifically for repeatability, also called test-retest reliability [28], the guarantees of a process discovery algorithm come into play. If the different samples are of similar quality, then the constructed models should be of similar quality. However, current process discovery algorithms do not explicitly claim to provide such guarantees, nor does an approach exist to study such claims. Such an approach can rely on the use of samples to study the relation between the input event log quality and the resulting process model quality.

In the next section, we propose to use sampling techniques to measure the quality of event logs, and show how sampling can be used to establish a relation between log quality and model quality in Section 4.

3 Event Log Sampling

As argued in the previous section, a necessary step in providing guarantees on the results of process discovery algorithms is to establish measures for log quality. In Section 3.1, we elaborate on the idea of considering an event log L to be a (non-random) sample of a continuous stream of events generated by a true process TP . As a consequence, any random sample of L is a sample of the same stream of events. Thus, random sampling techniques can be used to establish the quality of a sample event log with respect to L . Several random sampling techniques are presented in Section 3.2. Then, we propose seven requirements for measuring log quality in Section 3.3. Section 3.4 discusses to what extent standard error measures can be applied to quantify the lack of log quality. Last, we study the effect of random sampling techniques on log quality in Section 3.5.

3.1 Event Logs as Samples

We argue that any event log can be studied as a random sample of traces generated by the true process. Similar to [36], the true process can be represented as a set of traces with some trace likelihood function that assigns a probability to each trace. Consequently, any sample of an event log is again a sample of the true process, as proposed in [13]. We consider a sample log S of an event log

L to be a subset of the traces observed in the event log, i.e., $S(\sigma) \leq L(\sigma)$, for all traces $\sigma \in L$ and $S(\sigma) = 0$ if $\sigma \notin L$. This allows drawing different samples from a given event log, and then comparing these samples with the event log to analyze the quality of these samples. Little is known about the representativeness or quality of random samples in process mining [13, 41]. In the remainder of this section, we propose random sampling techniques to be used in process mining and provide measures to analyze the quality of a sample with respect to the original event log.

3.2 Sampling Techniques

In this section, we propose several sampling techniques that can be used to draw a sample from an event log, such that each trace in the event log has the same probability of being sampled. Consequently, samples obtained using these techniques can be used to estimate the characteristics of the event log and, thus, of the true process. The sampling techniques build on simple random sampling (Section 3.2.1) and stratified sampling (Section 3.2.2). An illustration of the discussed sampling techniques for the event log summarized in Table 1 is shown in Table 2. To obtain the samples, a sampling ratio of 25% is used.

Table 1: Example event log L with eight traces. The log consists of four unique traces.

Trace	$\langle a, d, g \rangle$	$\langle a, c, g \rangle$	$\langle a, b, g \rangle$	$\langle a, e, g \rangle$
Frequency	4	2	1	1

Table 2: Illustration of the different sampling techniques on the event log from Table 1, showing the challenges associated with handling infrequent traces. Each row represents an example sample log, given by the frequency of traces, constructed with the respective technique.

Sample	Technique	Sampled event log, sample ratio: 25%			
	Sequence Expected frequency	$\langle a, d, g \rangle$	$\langle a, c, g \rangle$	$\langle a, b, g \rangle$	$\langle a, e, g \rangle$
S_1	Random Fixed	1	0	1	0
S_2	Random Probability	2	1	1	0
S_3	Stratified	1	0	0	0
S_4	Existential Stratified	1	1	1	1
S_5	Stratified Plus	1	0	0	1
S_6	Stratified Squared	1	1	0	0

3.2.1 Simple Random Sampling Techniques

The first two sampling techniques are based on *simple random sampling*, where a sample is created by randomly including traces with a predetermined sampling ratio. *Random fixed* sampling starts by calculating the size of the event log, and then determines the size of the sample log. The sample log is then created by

randomly drawing traces from the event log until the sample log has the desired size. To illustrate the technique, two traces, sample log S_2 was created by drawing $\langle a, d, g \rangle$ and $\langle a, b, g \rangle$ out of the eight cases from event log L in Table 2.

Another sampling technique is *random probability*, where each trace is included in the sample based on the inclusion probability. For example, creating a sample of 25% results in each trace having a probability of 25% to be included in the sample log. As an example, sample log S_2 was drawn using this technique in Table 2. It has four cases: two instances of trace $\langle a, d, g \rangle$, and traces $\langle a, c, g \rangle$ and $\langle a, b, g \rangle$ were both drawn once. As the example shows, a disadvantage of this technique is that the resulting size of the sample might differ from the intended size.

3.2.2 Stratified Sampling Techniques

Stratified sampling takes a different approach to creating random samples. Classical *stratified sampling* [8] divides the data into unique groups called strata. A simple random sample is taken from each stratum. For process discovery, these strata can be formed based on unique traces. In theory, this sampling technique would give more representative samples because of the stratification of unique traces. However, one has to be careful when applying stratified sampling: as only a natural number of traces can be added to a sample, a trace can only be added fully or not at all. This technique is illustrated in Table 2: there are four strata. In the first, 25% of four sequences are selected, i.e., a single trace. For each of the other strata, the number of elements to select is lower than 1, i.e., no traces are selected from the other strata. Hence, a problem occurs if a stratum contains fewer traces than expected to be sampled. One way to solve this is by rounding, e.g. using the half to even rule (cf. IEEE 754).

Another solution for unsampled strata is *existential stratified sampling*. Similar to the classical stratified approach, the half to even rule is used. However, after rounding, a trace from each unsampled stratum is added to the sample log. Although it ensures that the directly-follows relations of the sample log and the original event log are identical, the main disadvantage is that these strata are an overrepresentation in the sample. As shown in Table 2, the stratified sample S_4 is complemented by adding a single trace from the remaining strata. Existential stratified sampling shows a trade-off between existential completeness of directly-follows relations and the representativeness of the frequencies of directly-follows relations.

The *stratified plus* sampling method tries to find a balance between existential completeness and frequency representativeness by randomly sampling additional cases whose trace has not been included in the sample yet. It uses the number of traces that were expected to be sampled and the number of traces sampled by stratified sampling in order to determine how many additional traces should be sampled. In Table 2, the stratified sample, containing one trace, is complemented by adding one trace, randomly selected from the remaining traces.

A different approach is taken in the *stratified squared* sampling approach. It

extends the classical stratified approach by randomly sampling additional traces that have not been included in the sample yet, based on the number of cases that were expected to be sampled and the number of cases sampled by stratified sampling: from the strata that are not represented, traces are randomly selected, until the sample log has the desired size. First, a stratified sample is drawn. Then, the number of sampled traces is compared to the number of expected traces based on the sampling ratio. Due to rounding, the number of expected traces can be greater than the number of actually added traces. If this happens, the uncovered strata are sorted based on their frequency, and a trace of each of these strata is added. This procedure continues until the number of sampled traces matches the expected number of traces or until all strata are covered.

3.3 Requirements for Sample Quality Measures for Process Mining

All random sampling techniques discussed in the previous section draw traces from the event log: for each trace it is decided whether the whole trace is added to the sample. Different approaches exist to estimate the quality of a sample, e.g., grounded in the Observed Trace variants Ratio (OTR) [19], i.e., the ratio of observed unique traces in the sample with respect to all unique traces in the original event log. As Table 2 shows, even though samples S_1 and S_6 both contain two out of four traces, the amount of information they contain is different, as S_6 contains the more frequent trace $\langle a, c, g \rangle$.

Most discovery algorithms (cf. [35, 37, 40, 14]) abstract from traces by using the directly-follows relation. Therefore, we propose, similar to [13, 3], to measure sample quality based on the directly-follows relation. The directly-follows relation $>_L$ is defined on pairs of events a and b , such that $a >_L b$ iff the event log L contains a trace in which the two activities a and b occur consecutively.

One way for comparing a sample to the original event log is *existential completeness*, i.e., the extent to which all possible directly-follows pairs are present in the sample. This results in the first sample quality measure: *coverage*. Coverage is defined as the ratio of unique directly-follows pairs present in the sample to the number of unique directly-follows pairs in the event log.

Coverage does not take the occurrence frequency of behavior into account. For example, sample logs S_1 , S_5 and S_6 all have a coverage of 50%. though sample log S_6 contains a more frequent trace than the other two samples. Different requirements can be defined to consider frequency representativeness in measuring sample quality. However, measuring the frequency representativeness of a sample is more subjective than measuring existential completeness. For example, for process conformance testing, like audit, rare behavior might be of interest, while for another type of project, only the most frequent traces are essential. Therefore, instead of pointing towards a single best measure for frequency representativeness, we present a list of generic requirements, and propose several measures, assessing them against these generic requirements. The

proposed requirements are formulated in terms of a penalty, or an *error*, that measures of sample quality should assign to samples under different conditions.

- R1. **Respect exact matches:** The measure should report no error when the frequencies of directly-follows pairs of the sample exactly match the expected frequencies;
- R2. **Doubling has no effect:** Doubling the number of unique directly-follows pairs present in the original event log should not affect the reported error when the new unique directly-follows pairs are equally often expected and sampled as the unique directly-follows pairs before doubling;
- R3. **Be proportional:** Doubling the number of occurrences of every directly-follows pair present in the original event log should not affect the reported error when the deviation of each sampled directly-follows pair is proportionally the same (e.g. the deviation of a directly-follows pair which is expected to occur five times, but is sampled three times is proportional to the same directly-follows pair being expected to occur fifty times, but being sampled thirty times);
- R4. **Punish absolute deviations:** When the sample size is varied while the absolute deviation is kept the same (e.g. all directly-follows pairs are off by one occurrence), then the error reported by the measure should increase when the sample size decreases;
- R5. **Punish large over small errors:** When one directly-follows pair is oversampled by four (i.e., is sampled four more times than its expected frequency), then the reported error should generally be larger compared to when four directly-follows pairs are oversampled by one occurrence;
- R6. **Trace frequency:** A sample where only the least often occurring directly-follows pair is off by one (i.e., sampled once more or once less often than its expected frequency) should generally report a higher error than the same sample where only the most often occurring directly-follows pair is off by one;
- R7. **Maintain perfect sampled pairs:** Given two samples of different sizes, if the frequency of a directly-follows pair matches the expected occurrences in both samples, and all other pairs have the same deviations, then the smaller sample should be penalized more.

3.4 Sample Quality Measures for Process Mining

In this section, we show how standard error measures from the field of statistics can be adapted to quantify log quality. Error measures are used to quantify the error between the expected values and the real occurrences. As argued in the previous section, we propose to measure sample quality based on the directly-follows relation. Thus, we want to establish the error based on the

directly-follows relation, and compare the occurrence of each behavior, i.e., each element in the directly-follows relation, with its expected value. Each expected value can be derived from the sampling ratio. As a result, we obtain several measures of sample quality. In the definitions that follow, by \mathbf{e} we denote the expected behavior, and by \mathbf{s} , we denote the sampled behavior as vectors of length n (i.e. n denotes the number of unique directly-follows pairs):

The (Normalised) Mean Absolute Error (NMAE) calculates the normalized absolute deviation (i.e., error) of the number of occurrences of each unique directly-follows relation of the sample from their respective expected frequency:

$$\text{NMAE} = \frac{\text{MAE}}{\text{avg } \mathbf{e}} = \frac{\sum_{i=1}^n |\mathbf{s}_i - \mathbf{e}_i|}{\sum_{i=1}^n \mathbf{e}_i}. \quad (1)$$

The (Normalised) Root Mean Square Error (NRMSE) is similar to the NMAE, but uses the root of the squared deviations, instead of the absolute values, thus penalizing large deviations more heavily:

$$\text{NRMSE} = \frac{\text{RMSE}}{\text{avg } \mathbf{e}} = \frac{\sqrt{\frac{1}{n} \sum_{i=1}^n (\mathbf{s}_i - \mathbf{e}_i)^2}}{\frac{1}{n} \sum_{i=1}^n \mathbf{e}_i}. \quad (2)$$

The Mean Absolute Percentage Error (MAPE) expresses the deviation as a percentage. Its symmetric version (sMAPE) has the advantage that the undersampling of behavior is penalized more heavily:

$$\text{MAPE} = \frac{1}{n} \sum_{i=1}^n \left| \frac{\mathbf{e}_i - \mathbf{s}_i}{\mathbf{e}_i} \right|, \quad \text{sMAPE} = \frac{1}{n} \sum_{i=1}^n \frac{|\mathbf{e}_i - \mathbf{s}_i|}{\mathbf{e}_i + \mathbf{s}_i}. \quad (3)$$

The Symmetric Root Mean Square Percentage Error (sRMSPE) is similar to sMAPE but uses the root mean square error instead of the mean absolute error, thus penalizing large deviations more heavily:

$$\text{sRMSPE} = \sqrt{\frac{1}{n} \sum_{i=1}^n \left(\frac{\mathbf{e}_i - \mathbf{s}_i}{\mathbf{e}_i + \mathbf{s}_i} \right)^2}. \quad (4)$$

These measures assess the behavioral quality of a sample with respect to the event log it is drawn from. In other words, these measures provide ways to establish the quality of the input of process discovery algorithms. Table 4 shows each measure on the samples in Table 2. The frequencies of behavior in both samples are shown in Table 3. Sample S_4 has perfect coverage, as every directly-follow pair of L occurs at least once in the sample. The MAE and the NMAE report the same relative error for S_1 and S_4 , as the expected frequencies are equal for both samples. Note that the NMAE would adjust itself

Table 3: The expected frequencies of directly-follows pairs together with the frequencies of directly-follows pairs of sample S_1 and sample S_4 of event log L (Table 2).

	Frequency						
	Expected	S_1	S_2	S_3	S_4	S_5	S_6
$a >_L b$	0.25	0	1	0	1	0	0
$a >_L c$	0.50	1	1	0	1	0	1
$a >_L d$	1.00	1	2	1	1	1	1
$a >_L e$	0.25	0	0	0	1	1	0
$b >_L g$	0.25	0	1	0	1	0	0
$c >_L g$	0.50	1	1	0	1	0	1
$d >_L g$	1.00	1	2	1	1	1	1
$e >_L g$	0.25	0	0	0	1	1	0

Table 4: Errors reported by the proposed measures on the samples in Table 2.

Error Measure	S_1	S_2	S_3	S_4	S_5	S_6
Coverage	0.50	0.75	0.25	1.00	0.50	0.50
MAE	0.25	0.63	0.25	0.50	0.38	0.25
NMAE	0.50	1.25	0.50	1.00	0.75	0.50
RMSE	0.31	0.68	0.31	0.59	0.47	0.31
NRMSE	0.61	1.37	0.61	1.17	0.94	0.61
MAPE	0.75	1.50	0.75	1.75	1.25	0.75
sMAPE	0.58	0.57	0.75	0.38	0.65	0.58
sRMSPE	0.73	0.63	0.87	0.46	0.77	0.73

with respect to sample size, whereas MAE is size agnostic. Sample S_4 scores better on sMAPE than S_1 , as S_1 does not sample two directly-follows pairs, whereas S_4 only oversamples pairs. This illustrates that the sMAPE measure gives a higher penalty for unsampled behavior. RMSE, NRMSE, and sRMSPE give comparable results for these two samples as these samples do not contain large deviations between actual and expected frequencies.

As these examples show, not all measures satisfy all requirements. Therefore, we analyzed each measure against the requirements. The results of this analysis is shown in Table 5. A shortcoming of the MAE measure is that changes in the expected sample size are not reflected in the reported error (R4). For example, in one sample, a directly-follows relation is expected to occur ten times and occurs nine times, while in another sample with larger sample size, this directly-follows relation is expected to occur one hundred times and occurs ninety-nine times. The MAE gives these two samples an equal error because both are exactly off by one. It fails to satisfy most requirements, see Table 5. Normalizing the MEA, i.e., the NMAE measure, results in a measure that leads to the fulfillment of most of the requirements.

The RMSE measure behaves similarly to MAE. However, it penalizes more

Table 5: Testing each frequency representativeness measure against the requirements.

Measure	R1	R2	R3	R4	R5	R6	R7
MAE	✓	✓	✗	✗	✗	✗	✗
NMAE	✓	✓	✓	✓	✗	✗	✓
RMSE	✓	✓	✗	✗	✓	✗	✗
NRMSE	✓	✓	✓	✓	✓	✗	✓
MAPE	✓	✓	✓	✓	✗	✓	✗
sMAPE	✓	✓	✓	✓	✗	✓	✗
sRMSPE	✓	✓	✓	✓	✓	✓	✗

significant deviations more heavily, which can be a desired property if unbalanced samples are undesired (i.e. samples where the number of occurrences of one or a few directly-follows relations deviate a lot from their expected frequency). Its normalization, i.e., the NRMSE, results in a measure that satisfies all requirements, except R6.

The main difference between the MAPE and NMAE is that the MAPE does not decrease the error when increasing the number of occurrences of one or more perfectly sampled directly-follows relations while still keeping them perfectly sampled. Vice versa, the NMAE does not report a lower error when the most occurring directly-follows relation is off by one compared to the least occurring directly-follows relation being off by one. Its symmetric version, i.e., the sMAPE measure, ticks the same requirements, but the sMAPE does favor existential completeness compared to the MAPE, because it gives unsampled behavior the highest possible penalty. This makes the sMAPE measure more appropriate for process discovery of rare behavior.

Comparing the sRMSPE with the sMAPE shows that the former uses the square root error, instead of the mean absolute error. Consequently, the sRMSPE gives a more significant penalty to directly-follows relations whose number of occurrences is further off its expected frequency. If this property is desired, then the sRMSPE should be selected over the sMAPE.

Overall, the analysis of the requirements shows that if existential completeness is important, sMAPE or sRMSPE should be chosen; otherwise, the NMAE or NRMSE should be used. When single large deviations are not desired, then the root mean square error-based measures should be used instead of the mean absolute error variants.

3.5 Effects of Sampling Techniques on Log Quality

Using the sample quality measures, the effects of the different sampling methods can be studied. For this evaluation, we used two event logs that were generated from a Petri net with a start transition a , followed by five parallel transitions b , c , d , and e , and final transition f . The first event log, L_1 , is a balanced log, i.e., all traces have a similar frequency, whereas the second event log, L_2 , has

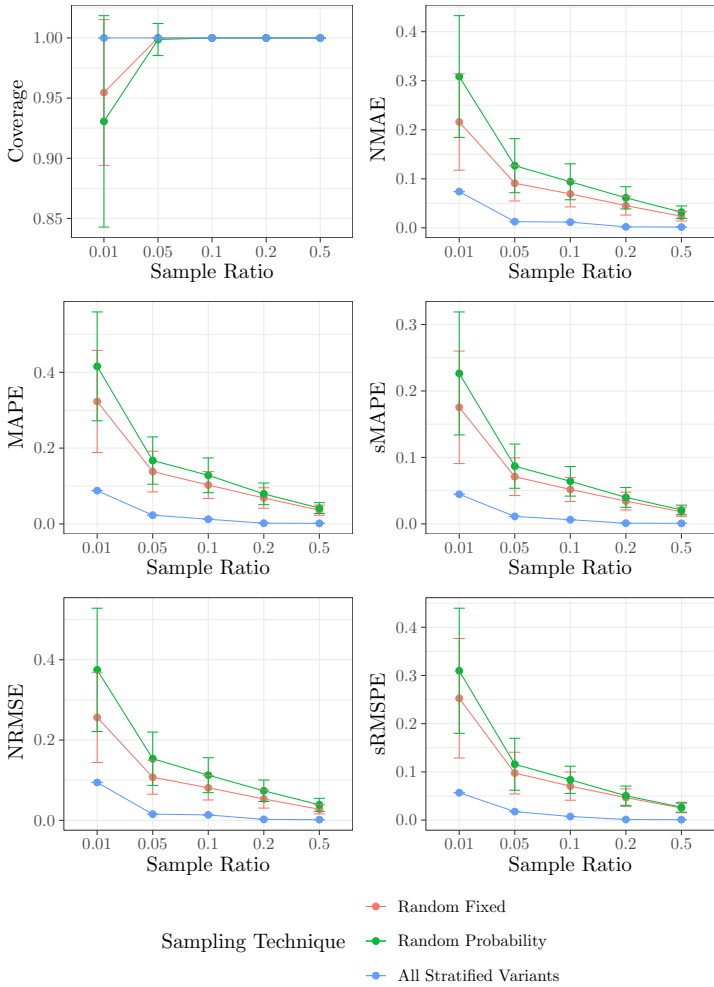


Figure 2: The effects of different sample ratios and sampling techniques on the quality measures of samples from the balanced event log L_1

many infrequent traces, i.e., all traces occur only once in the event log.

Each event log has been sampled using these random sampling techniques: random sampling with a fixed sample size (random fixed), probability-based random sampling (random probability), stratified sampling, existential stratified sampling, stratified plus sampling, and stratified squared sampling. Sampling with each technique was repeated one hundred times for each of the following five sampling ratios: 0.01, 0.05, 0.1, 0.2, and 0.5. This resulted in one hundred samples for each combination of sampling technique and sample ratio. For each sample, the coverage, NMAE, MAPE, sMAPE, NRMSE, and sRMSPE have

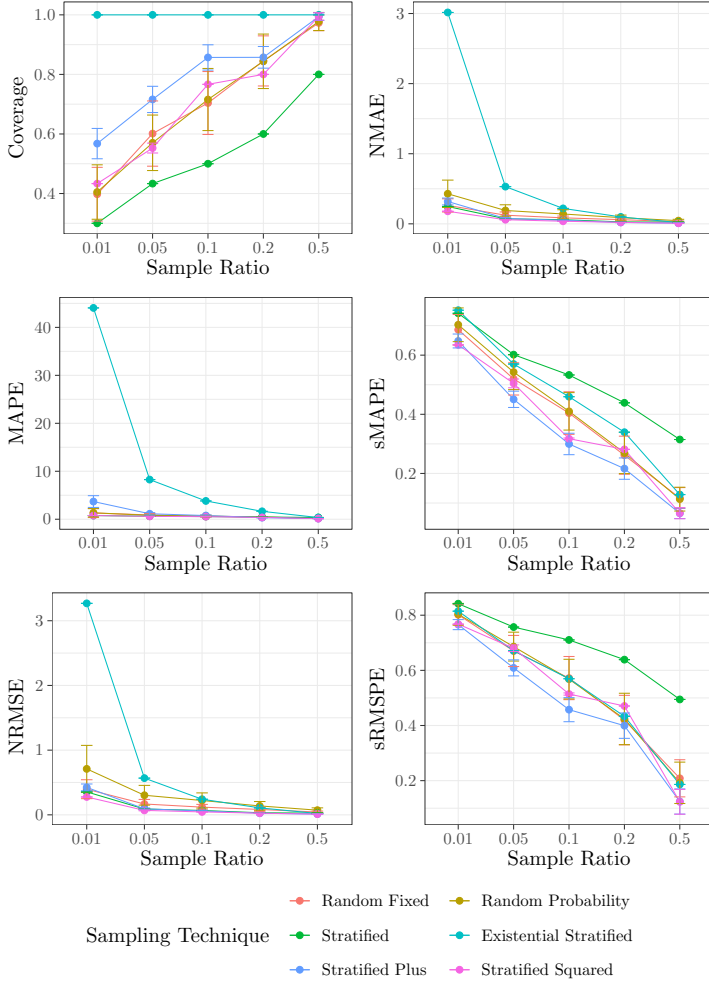


Figure 3: The effects of different sample ratios and sampling techniques on the quality measures of samples from event log L_2 containing only infrequent traces.

been calculated. Next, for each combination of sampling technique and sample ratio, the quality measures have been averaged over the one hundred samples and the standard deviation has been calculated.

Figure 2 shows the effects of sampling on the balanced event log. As there are no infrequent traces in this event log, the values reported by all four different variations of stratified sampling techniques are exactly the same. The probability-based random sampling technique consistently performs worst on all measures. The fixed sample size random sampling technique only performs marginally better. All four stratified sampling techniques seem to create near-

perfect samples, especially when the sample ratio is 0.05 or larger.

The effects of the sampling techniques are more evident in the event log with infrequent traces, as shown in Figure 3. By definition, existential stratified sampling always results in a perfect coverage of 1. However, as the measures show, it oversamples rare directly-follows relations, which is especially true for the smaller sample ratios. Stratified sampling performs the worst of all sampling techniques, as it leaves out all rare sequences from the original event log. The evaluation also confirms that the non-symmetric measures (NMAE, MAPE, and NRMSE) perform worse on an event log with many infrequent traces than the symmetric measures (sMAPE, sRMSPE). Probability-based random sampling performs poorly on these sample measures, while stratified squared sampling consistently has the lowest error on NMAE and NRMSE. The difference between the sampling techniques is most significant with a sample ratio of 0.01, while for larger sample ratios, the difference between the sampling techniques decreases.

As this evaluation shows, random sampling and error measures can be used to express log quality, given an original event log. In the next section, we propose four maturity stages for process discovery algorithms, and show how sampling can be used to establish a relationship between log and model quality.

4 Designing Process Discovery Algorithms with Guarantees

As observed in a study on the quality of conformance measures [23], some process discovery algorithms have a large variability in the quality of the constructed process models. In particular, given different samples of a single event log, the same algorithm sometimes provides good results on small samples, while on larger samples, the algorithm discovers worse models. On further inspection, these algorithms are state of the art and do not perform any major “process mining crimes” [25]. In addition, they “glitter” in the benchmark study reported in [2].

We consider this observation a threat to the application of process mining, particularly for its repeatability and, hence, the reliability of its results. Suppose for a true process, several event logs are captured and analyzed, and the results do not agree – they differ in quality. Several explanations for this phenomenon are possible. A first explanation could be the quality of the input, i.e., the quality of the event logs differed significantly. However, as the observation highlights, another plausible – yet undesirable – explanation lies in the process discovery algorithm itself. In other words, if the process discovery algorithm does not provide any guarantees on the quality of the resulting models, it is impossible to exclude the algorithm as a root cause.

Consequently, we advocate that process discovery algorithms should provide guarantees on the quality of the produced results. To this end, we propose to distinguish four stages of maturity for a process discovery algorithm:

1. The algorithm is well designed;

2. The algorithm is validated on real-life input;
3. The algorithm has an established relationship between the log and model quality;
4. The algorithm is effective.

As illustrated later, not all algorithms make it to the second stage. Arguably, algorithms that do not pass the second stage should not be used in empirical studies. The third and fourth stages have never been considered before in the context of process discovery. Once the algorithm is shown to be applicable on real-life examples while producing useful results, the authors should study which guarantees their algorithm provides in a controlled setting where the true process is known. To satisfy the requirements of the last stage, the algorithm should provide evidence that in settings where the true process is unknown, the algorithm provides the guarantees stated in stage 3.

4.1 Stage 1: The Algorithm is Well Designed

In the first stage, the developers of a process discovery algorithm should properly introduce their algorithm, by providing the following:

- Criteria on the event logs the algorithm uses as input, e.g., requirements on the true process that generates the event logs;
- The class of process models the algorithm constructs;
- Evidence for meeting the quality goals of the algorithm;
- An initial evaluation of the algorithm on artificial data sets.

Most process discovery algorithms satisfy the requirements of this stage. For example, the ILP-miner [37] is designed for the class of classical Petri nets with interleaving semantics. It meets its quality goals: it is demonstrated that the ILP miner always returns a Petri net with a perfect recall score for any input. It imposes no requirements on the input event logs and is tested on artificial logs. Also, the α -miner algorithm [35] is at least in this stage. It is designed with the rediscoverability of well-structured Workflow nets as a goal. To rediscover the true process, it imposes two requirements on an input event log: it should contain all directly-follows relations present in the true process, and the true process should be block-structured [20]. A similar argument holds for the Inductive Miner [15].

4.2 Stage 2: The Algorithm is Validated

Even though an algorithm may be well designed, and passes stage 1, it is not guaranteed that it is useful in practice. Therefore, the second stage of process discovery algorithm maturity is concerned with the validation of the algorithm

on a collection of real-life event logs, such as logs used in the benchmark reported in [2]. Several existing algorithms fail to reach this stage. For example, the α -miner is, theoretically, a robust algorithm, but the requirements it imposes on the true process are often too strong for application in real-life situations [34]. Similarly, the ILP-miner, despite being theoretically grounded, has limitations for when it comes to practice, primarily because of its guaranteed recall and runtime performance. Other algorithms, such as the Inductive Miner [15], the Declare Miner [17], and the Split Miner [1] have been applied successfully on several real-life event logs, and, thus, pass this stage.

4.3 Stage 3: The algorithm has an Established Relationship Between Log and Model Quality

Although passing the second stage shows the algorithm’s capabilities, this provides little guarantee on the quality of the constructed process models, in general, for a wide range of input event logs. As a first step in establishing a relationship between the input event log and the output model qualities, one can demonstrate to what degree the algorithm supports the principles sketched in Figure 1. In other words, authors of the algorithm need to show that if an event log is a faithful representation of the true process, as per measure \mathcal{P}^T , then the algorithm satisfies properties similar to those listed below:

- P1. For a sample log S that approaches the perfect quality, the quality \mathcal{P}^S of the model discovered from S approaches \mathcal{P}^T ;
- P2. For two sample logs S_1 and S_2 , if S_1 has a better quality than S_2 , then the model quality \mathcal{P}^{S_1} of the model discovered from S_1 is better than the quality \mathcal{P}^{S_2} of the model discovered from S_2 .

Authors of a process discovery algorithm can follow different strategies to provide evidence for these properties. The most potent form of evidence is formal proof that the algorithm satisfies these properties for specific instantiations of log and model quality measures. That way, a strong relationship between

Algorithm 1: Establish Relation

```

1 while True do
2   TP ← GenerateTrueProcess( $\mathcal{M}$ , A);
3   foreach  $i \in [1..N]$  do
4     L ← GenerateLog(TP);
5      $\mathcal{P}^T \leftarrow \text{computeModelQuality}(L, TP)$ ;
6     foreach  $r \in \text{ratios}$  do
7       foreach  $j \in [1..K]$  do
8         S ← DrawSample(L, r);
9         e ← computeSampleQuality(L, S);
10        M ← DiscoverModel(S);
11         $\mathcal{P}^S \leftarrow \text{computeModelQuality}(S, M)$ ;

```

Table 6: Results of the controlled experiment, showing the Spearman rank correlation between the error measures and precision. All bold values are statistically significant ($p < 0.001$).

Model	True Process		Precision			
	precision	Cov.	sMAPE	sRMSPE	NRMSE	NMAE
1	0.538	0.658	-0.988	-0.986	-0.988	-0.989
2	0.797	0.470	-0.986	-0.985	-0.901	-0.954
3	0.935	0.781	-0.990	-0.989	-0.975	-0.984
4	0.953	0.705	-0.991	-0.992	-0.984	-0.987
5	0.988	0.540	-0.983	-0.981	-0.980	-0.986
6	0.871	0.532	-0.934	-0.938	-0.917	-0.926
7	0.943	0.511	-0.991	-0.989	-0.986	-0.989
8	0.616	0.773	-0.992	-0.991	-0.989	-0.990
9	0.710	0.519	-0.981	-0.978	-0.970	-0.973
10	0.883	0.703	-0.982	-0.982	-0.977	-0.976

Table 7: Results of the controlled experiment, showing the Spearman rank correlation between the error measures and recall. All bold values are statistically significant ($p < 0.001$).

Model	True Process		Recall			
	recall	Cov.	sMAPE	sRMSPE	NRMSE	NMAE
1	1.000	0.338	-0.356	-0.354	-0.354	-0.356
2	1.000	0.154	-0.051	-0.052	0.012	-0.004
3	1.000	0.637	-0.406	-0.417	-0.410	-0.412
4	1.000	-0.103	0.105	0.108	0.081	0.090
5	1.000	0.437	-0.201	-0.206	-0.207	-0.201
6	1.000	-0.529	0.973	0.962	0.963	0.968
7	1.000	0.456	-0.242	-0.240	-0.228	-0.231
8	1.000	0.114	-0.148	-0.154	-0.156	-0.157
9	1.000	0.518	-0.327	-0.330	-0.340	-0.341
10	1.000	0.116	-0.022	-0.027	-0.016	-0.023

an input log quality and the resulting model quality can be established. We also encourage algorithm designers to define algorithm-specific log quality measures. If formal proof is not feasible, statistical evidence of these properties can be provided. To this end, we propose a controlled experiment as outlined in Algorithm 1. Such controlled experiment implements the approach outlined in Figure 1. It requires a model generator for the class of true processes \mathcal{M} the algorithm supports and a set of activities A . The algorithm then generates repeatedly for a true process one or more event logs, and for each event log a set of samples to relate log and model quality.

We propose to use statistical tests to evaluate the two properties. Property P1 calls for analyzing the relationship between the expected \mathcal{P}^T and the observed \mathcal{P}^S . To establish property P2, for instance, the Spearman rank correlation can be used to test whether there is a strong correlation between the sample quality and the model quality. If this is the case, then statistical evidence has been provided for the relationship between log and model quality.

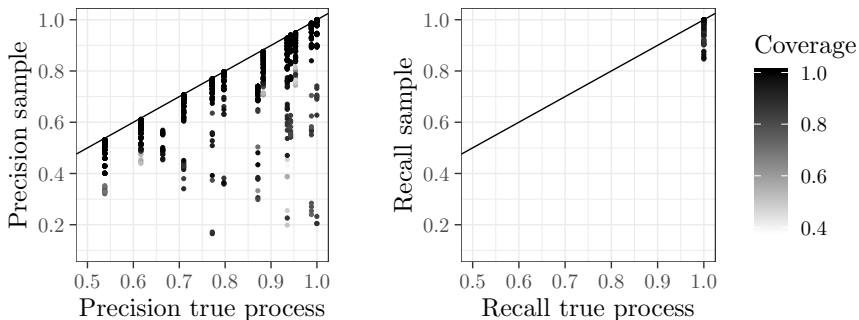


Figure 4: Relation between the quality of the true process and the quality of the discovered models, for precision (left) and recall (right). Darker points represent a higher coverage.

4.3.1 Example Evaluation

To show the feasibility of the approach, the controlled experiment has been implemented in ProM¹ for the Inductive Miner [15]. Precision and recall are calculated using an implementation of exact matching entropy-based measures in Entropia [21]. For each true process, a single event log with 5,000 traces has been generated. The event logs were 10 times sampled for 12 sampling ratios: 0.01, 0.02, 0.05, and 0.1 up to 0.9.

The results are shown in Tables 6 and 7, and in Figure 4. The diagonal lines in the graphs indicate the desired result: a point on the diagonal indicates that both the expected value and the actual value coincide. The darkness of each data point indicates the coverage: the darker the point is, the higher the coverage. As the results show, both for precision and recall, most values are on this diagonal. Therefore, we can conclude that property P1 holds for precision and recall.

For each model that describes the true process, the Spearman rank correlation is calculated between each of the log quality measures and precision, and similarly for recall. As for the measures sMAPE, sRMSPE, NRMSE, and NMAE, 0 is the best quality, a negative correlation indicates the required guarantee that samples of higher quality result in better discovered models, whereas for coverage, a positive correlation indicates this result. As can be seen in Tables 6 and 7, the experiment generates mixed results. Though property P2 holds for precision, it is not satisfied for recall. Hence, we can conclude that the Inductive Miner satisfies the two properties for precision, but fails to do so for recall on the second property.

¹The source code is available on: <https://github.com/ArchitectureMining/SamplingFramework>

4.4 Stage 4: The Algorithm is Effective

An established relationship between log and model quality, the essence of the third stage, does not guarantee the algorithm to be effective in real-life situations. The main caveat in the controlled environment of the previous stage is that the true process is known. Each event log is generated from the known true processes. In real-life situations, the true process is unknown, and, hence, may invalidate assumptions of the discovery algorithm. For example, the Inductive Miner assumes event logs to be generated from process trees. However, no criteria are given to test whether an event log is generated by a process tree, nor does the algorithm provide any details on the model quality if the assumption is invalid.

In this stage, the algorithm designer has to validate how effective the algorithm is in real-life situations. One way to obtain insights into the effectiveness of the algorithm is to apply sampling on a benchmark. This benchmark can be a set of well-known real-life event logs as used in [2], or can be generated automatically, if the designers ensure that the class of generated models is larger than the class of true processes studied in the previous stage. The algorithm designers need to analyze property P2 in the absence of a true process. In other words, even if the true process is unknown, event logs of better quality should return better quality models. This may result in an experiment as outlined in Algorithm 2.

The analysis of property P2 in the absence of a true process can have two possible outcomes. Either it is shown that the algorithm has the desired property, or, if this is not possible, the algorithm should be further improved, or provide additional log quality measures, that guarantee that an event log satisfies the assumptions of the process discovery algorithm.

4.4.1 Example Evaluation

As an example of the analysis in stage 4, we conducted the proposed experiment on the Inductive Miner [15]. Two real-life event logs have been selected, the Road Traffic Fine event log [10] and the Sepsis event log [18]. The Road Traffic Fine log has in total 150,370 traces and 561,470 events. There are 231 unique traces and 11 unique event types. The Sepsis log consists of 1,049 traces, of which 845 are unique, and 15,190 events with 16 unique event types. Sampling

Algorithm 2: Test Effectiveness

```
1 foreach  $L \in \text{Benchmark}$  do
2   foreach  $r \in \text{ratios}$  do
3     foreach  $j \in [1..K]$  do
4        $S \leftarrow \text{DrawSample}(L, r)$ ;
5        $e \leftarrow \text{computeSampleQuality}(L, S)$ ;
6        $M \leftarrow \text{DiscoverModel}(S)$ ;
7        $\mathcal{P}^S \leftarrow \text{computeModelQuality}(S, M)$ ;
```

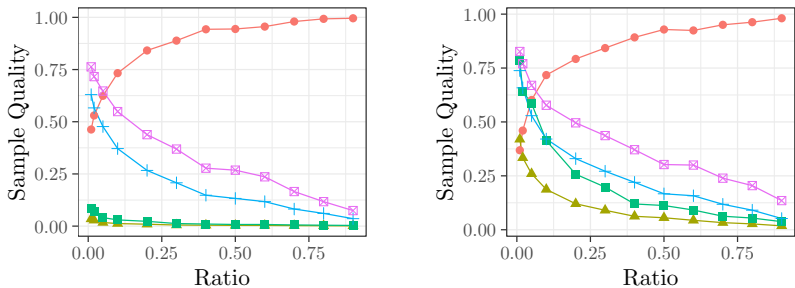


Figure 5: Plot of ratio and the sample quality measures coverage (●), sMAPE (+), sRMSPE (⊠), NRMSE (■) and NMAE (▲) for the Road Traffic Fine log (left) and the Sepsis log (right).

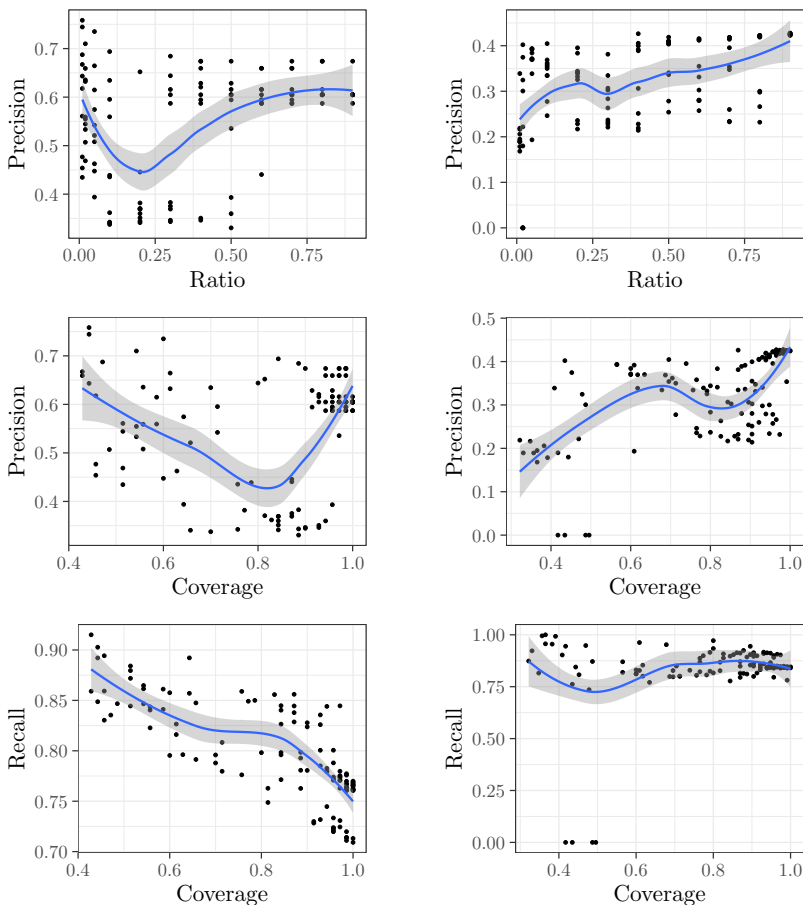


Figure 6: Plots of ratio and precision, and coverage with precision and recall for the Road Traffic Fine log (left) and the Sepsis Log (right).

was done at the same sampling ratios as before: 0.01, 0.02, 0.05, and 0.1 up to 0.9. For each ratio, ten samples were drawn.

The sample quality measures for the Road Traffic Fine log are shown on the left in Fig. 5. As the plot shows, the larger the sampling ratio, and thus the log size, the better the quality is (error measures: $\rho < -0.9$, $p < 0.001$, coverage: $\rho = 0.96$, $p < 0.001$). Sample size and the conformance measure on precision (Fig. 6) show a moderate positive correlation ($\rho = 0.56$, $p < 0.001$), while there is no correlation between sampling ratio and recall ($\rho = 0.03$, $p = 0.72$). Analyzing the quality measures with the conformance measures shows a different story. In Fig. 6, the coverage is plotted against the precision, indicating there is no correlation between coverage and precision. Further analysis revealed no correlations between the sample quality measures and precision (sMAPE: $\rho = -0.19$, $p = 0.03$, sRMSPE: $\rho = -0.18$, $p = 0.051$, NRMSE: $\rho = -0.21$, $p = 0.02$, NMAE: $\rho = -0.20$, $p = 0.03$, coverage: $\rho = 0.17$, $p = 0.06$). The correlations found for recall show that samples of worse quality result in better models (sMAPE: $\rho = 0.80$, $p < 0.001$, sRMSPE: $\rho = 0.79$, $p < 0.001$, NRMSE: $\rho = 0.77$, $p < 0.001$, NMAE: $\rho = 0.78$, $p < 0.001$, coverage: $\rho = -0.79$, $p < 0.001$).

For the Sepsis log, similar results are found. As indicated by the plots at the right hand side of Fig. 5, a correlation is found between the sampling ratio and the log quality measures (for all error measures: $\rho < -0.9$, $p < 0.001$, coverage: $\rho = 0.59$, $p < 0.001$). The larger the sampling ratio, the higher the precision is ($\rho = 0.57$, $p < 0.001$), but no correlation was found between sampling ratio and recall ($\rho = 0.03$, $p = 0.72$). A moderate negative correlation was found between the log quality measures and precision (for the error measures: $-0.60 < \rho < -0.50$, $p < 0.001$, coverage: $\rho = 0.59$, $p < 0.001$), while the log quality measures did not show any correlation with recall (for all measures: $-0.04 < \rho < 0.02$, $p > 0.70$).

As the results suggest, there is no clear relation between log and model quality. Hence, it is with the current measures not possible to conclude that the Inductive Miner is guaranteed to be effective in real-life situations. As a next step, new log quality measures should be developed that do establish the required relationship between log and model quality. The process can then be repeated until sufficient guarantees can be provided on the effectiveness of the algorithm.

5 Related Work

The statistical approach we propose to establish a relation between log and model quality relates to event data quality in general, builds upon established properties of conformance measures, and requires sampling techniques on event logs. This section reviews literature on these topics, and shows how our approach relates to them.

Measuring log quality. As the process mining manifesto articulates, process mining treats data as first-class citizens [33], and defines four data qualities, of which *completeness* is studied most. For example, [5] identifies four categories of process characteristics and 27 classes of event log quality issues. Most studies on event log quality focus on the incompleteness of the data. Examples include not having enough information recorded in the event log (e.g., missing cases or events) [5, 32], not having recorded enough behavior in the event log [12], or the traces not being representative of the process [12], and noise. Different notions of noise are studied, such as infrequent behavior that is either incorrect or rare [11]. However, event logs are studied in isolation in these studies. Instead, we argue to assess the quality of event logs relative to other event logs, using statistical techniques based on sampling.

Properties of conformance measures. The process mining community has recently initiated a discussion on which formal properties should “good” conformance measures satisfy. In [30], the authors proposed five properties for precision measures. For instance, one property states that for two process models that describe all the traces in the log, a less permissive model should not be qualified as less precise. By demonstrating that a measure fulfills such properties, one establishes its usefulness. In [23], the authors strengthened the properties from [30]. For example, according to these properties, the less permissive model from the example above should be classified as more precise. In [36], the precision properties from [30] were refined, and further desired properties for recall and generalization measures were introduced, resulting in 21 conformance propositions. Finally, in [22], properties for precision and recall measures that account for the partial matching of traces, i.e., traces that are not the same but share some subsequences of activities, were introduced. The precision and recall measures used in our evaluations satisfy all the introduced desired properties for the corresponding measures [30, 23, 36, 29].

Sampling in process mining. Sampling has been studied before in process mining, but never as a systematic approach to evaluate process discovery techniques. A first set of measures for the representativeness of samples have been proposed in [13]. Their results show the need for a systematic approach as proposed in this paper.

In [4], a sampling technique specific for the Heuristics Miner is described, claiming that only 3% of the original log is sufficient to discover 95% of the dependency relations. However, a proper evaluation of this claim has not been provided, nor are the results generalizable to other process discovery techniques.

A statistical framework based on *information saturation* is proposed in [3]. Their approach differs from the probability sampling techniques we propose. Instead of generating samples that estimate the event log, their approach focuses on creating a sufficiently small sample that contains as much information from the event log as possible. Consequently, this approach cannot be used to measure sample quality with respect to the event log.

Several biased sampling techniques are described in [26]. These techniques have been evaluated on six real-life event logs and three discovery techniques. The evaluation showed that sampling sometimes improves the F-measure for some of the models. A similar result on the F-measure was obtained in [16]. Their study applied the Google PageRank algorithm on event logs to create a representative sample, which reduced the execution time of the Inductive Miner by half without decreasing the F-measure. As the F-measure harmonizes precision and recall, and no analysis was performed on the reasons behind the improvements, it is unclear how sampling influenced the process discovery results of both studies. Instead of using sampling to improve the quality of the output, we propose to use probability sampling to analyze the input of algorithms, and to establish a relationship between log and model quality. This relationship then allows one to explore why some samples give better models than other samples.

6 Conclusion

This paper identifies the need for process discovery algorithms with guarantees that characterize the dependency between the quality of input event logs and the quality of the process models constructed from these event logs. In particular, we argue that process discovery algorithms should produce better models from better input logs. Currently, process discovery algorithms have never provided such guarantees, since, so far, we, as a community, lacked a theoretical foundation to establish such a relationship. In this paper, for the first time, measures for the statistical sample quality for ranking the quality of event logs are proposed. We recommend using grounded conformance checking measures for assessing the quality of the discovered models. Combining log quality measures with conformance measures provides a framework to formally define properties that express the desired guarantee that better event logs result in better models. These properties can be instantiated with various measures for quality of event logs and process models and be less or more pronounced, for example, imposing a strictly increasing or non-decreasing relation, or requiring a statistical association of a certain degree between the qualities of the corresponding logs and models. To overcome this problem, we propose four stages in the design of an algorithm. Each design comes with additional properties and obligations to establish effective algorithms with guarantees.

We invite the process mining community to further contribute to the discussion of desired qualities for process discovery algorithms to ensure that state-of-the-art algorithms fulfill them, and in this way, advance the field of process discovery as well as the design and evaluation of such algorithms.

Acknowledgement. Artem Polyvyanyy was in part supported by the Australian Research Council project DP220101516.

References

- [1] A. Augusto, R. Conforti, M. Dumas, and M. La Rosa. Split miner: Discovering accurate and simple business process models from event logs. In *ICDM 2017*, pages 1–10. IEEE, 2017.
- [2] A. Augusto, R. Conforti, M. Dumas, M. La Rosa, F. M. Maggi, A. Marrella, M. Mecella, and A. Soo. Automated discovery of process models from event logs: Review and benchmark. *IEEE Trans. Knowl. Data Eng.*, 31(4):686–705, 2019.
- [3] M. Bauer, A. Senderovich, A. Gal, L. Grunske, and M. Weidlich. How much event data is enough? a statistical framework for process discovery. In *CAiSE 2018*, volume 10816 of *LNCS*, pages 239–256. Springer, 2018.
- [4] A. Berti. Statistical sampling in process mining discovery. In *eKNOW 2017*, pages 41–43. IARIA, 2017.
- [5] J. C. Bose, R. S. Mans, and W. M. P. van der Aalst. Wanna improve process mining results? In *CIDM 2013*, pages 127–134. IEEE, 2013.
- [6] M. Bozkaya, J. M. A. M. Gabriels, and J. M. E. M. van der Werf. Process diagnostics : a method based on process mining. In *eKNOW 2009*, pages 22–27. IEEE, 2009.
- [7] J. C. A. M. Buijs, B. F. van Dongen, and W. M. P. van der Aalst. Quality dimensions in process discovery: The importance of fitness, precision, generalization and simplicity. *International Journal of Cooperative Information Systems*, 23(1), 2014.
- [8] William G Cochran. *Sampling techniques*. John Wiley & Sons, 1977.
- [9] Th. H. Cormen, Ch. E. Leiserson, R. L. Rivest, and C. Stein. *Introduction to Algorithms*. MIT Press Ltd, 2009.
- [10] M. de Leoni and F. Mannhardt. Road Traffic Fine Management Process, 2 2015. doi:10.4121/uuid:270fd440-1057-4fb9-89a9-b699b47990f5.
- [11] A. K. A. de Medeiros, A. J. M. M. Weijters, and W. M. P. van der Aalst. Genetic process mining: an experimental evaluation. *Data Min. Knowl. Discov.*, 14(2):245–304, 2007.
- [12] C. Günther. *Process mining in flexible environments*. PhD thesis, Eindhoven University of Technology, 2009.
- [13] B. Knols and J. M. E. M. van der Werf. Measuring the behavioral quality of log sampling. In *ICPM 2019*, pages 97–104. IEEE, 2019.
- [14] S. J. J. Leemans, D. Fahland, and W. M. P. van der Aalst. Discovering block-structured process models from event logs - A constructive approach. In *Petri Nets 2013*, volume 7927 of *LNCS*, pages 311–329. Springer, 2013.

- [15] S. J. J. Leemans, D. Fahland, and W. M. P. van der Aalst. Scalable process discovery with guarantees. In *EMMSAD 2015*, volume 214 of *LNBIP*, pages 85–101. Springer, 2015.
- [16] C. Liu, Y. Pei, Q. Zeng, and H. Duan. Logrank: An approach to sample business process event log for efficient discovery. In *Knowledge Science, Engineering and Management*, volume 11061 of *LNCS*, pages 415–425. Springer, 2018.
- [17] F. M. Maggi, J. C. Bose, and W. M. P. van der Aalst. Efficient discovery of understandable declarative process models from event logs. In *CAiSE 2012*, volume 7328 of *LNCS*, pages 270–285. Springer, 2012.
- [18] F. Mannhardt. Sepsis Cases - Event Log, 12 2016. doi:10.4121/uuid:915d2bfb-7e84-49ad-a286-dc35f063a460.
- [19] J. Pei, L. Wen, H. Yang, J. Wang, and X. Ye. Estimating global completeness of event logs: A comparative study. *IEEE Transactions on Services Computing*, 2018.
- [20] A. Polyvyanyy. *Structuring process models*. PhD thesis, University of Potsdam, 2012.
- [21] A. Polyvyanyy, H. Alkhamash, C. Di Ciccio, L. García-Bañuelos, A. A. Kalenkova, S. J. J. Leemans, J. Mendling, A. Moffat, and M. Weidlich. Entropia: A family of entropy-based conformance checking measures for process mining. In *ICPM Doctoral Consortium and Tool Demonstration*, volume 2703 of *CEUR*, pages 39–42. CEUR-WS.org, 2020.
- [22] A. Polyvyanyy and A. A. Kalenkova. Monotone conformance checking for partially matching designed and observed processes. In *ICPM 2019*, pages 81–88, 2019.
- [23] A. Polyvyanyy, A. Solti, M. Weidlich, C. Di Ciccio, and J. Mendling. Monotone precision and recall measures for comparing executions and specifications of dynamic systems. *ACM Trans. Softw. Eng. Methodol.*, 29(3):17:1–17:41, 2020.
- [24] Artem Polyvyanyy, Alistair Moffat, and Luciano García-Bañuelos. Bootstrapping generalization of process models discovered from event data. In *Advanced Information Systems Engineering 2022*, volume 13295 of *LNCS*, pages 36–54. Springer, 2022.
- [25] J. Rehse and P. Fettke. Process mining crimes - A threat to the validity of process discovery evaluations. In *BPM Forum 2018*, volume 329 of *LNBIP*, pages 3–19. Springer, 2018.
- [26] M. Fani Sani, S. J. van Zelst, and W. M. P. van der Aalst. Improving the performance of process discovery algorithms by instance selection. *Comput. Sci. Inf. Syst.*, 17(3):927–958, 2020.

- [27] William Shadish, Thomas Cook, and Donald Campbell. *Experimental and quasi-experimental designs for generalized causal inference*. Wadsworth Cengage Learning, 2002.
- [28] Peter Swanborn. A common base for quality control criteria in quantitative and qualitative research. *Quality & Quantity*, 30(1):19–35, 1996.
- [29] A. F. Syring, N. Tax, and W. M. P. van der Aalst. Evaluating conformance measures in process mining using conformance propositions. *ToPNOC*, pages 192–221, 2019.
- [30] N. Tax, X. Lu, N. Sidorova, D. Fahland, and W. M. P. van der Aalst. The imprecisions of precision measures in process mining. *Inf. Process. Lett.*, 135:1–8, 2018.
- [31] Andrei Tour, Artem Polyvyanyy, and Anna A. Kalenkova. Agent system mining: Vision, benefits, and challenges. *IEEE Access*, 9:99480–99494, 2021.
- [32] W. M. P. van der Aalst. *Process Mining—Data Science in Action, Second Edition*. Springer Berlin Heidelberg, 2016.
- [33] W. M. P. van der Aalst et al. Process mining manifesto. In *BPM Workshops*, volume 99 of *LNBIP*, pages 169–194. Springer, 2011.
- [34] W. M. P. van der Aalst, B. F. van Dongen, J. Herbst, L. Maruster, G. Schim, and A. J. M. M. Weijters. Workflow mining: a survey of issues and approaches. *Data and Knowledge Engineering*, 47(2):237–267, 2003.
- [35] W. M. P. van der Aalst, A. J. M. M. Weijters, and L. Maruster. Workflow Mining: Discovering Process Models from Event Logs. *Knowledge & Data Engineering*, 16(9):1128–1142, 2004.
- [36] Wil M. P. van der Aalst. Relating process models and event logs—21 conformance propositions. In *ATAED*, volume 2115 of *CEUR Workshop Proceedings*, pages 56–74. CEUR-WS.org, 2018.
- [37] J. M. E. M. van der Werf, B. F. van Dongen, C. A. J. Hurkens, and A. Serebrenik. Process discovery using integer linear programming. *Fundamenta Informaticae*, 94(3-4):387 – 412, 2009.
- [38] Jan Martijn E. M. van der Werf, Artem Polyvyanyy, Bart R. van Wensveen, Matthieu J. S. Brinkhuis, and Hajo A. Reijers. All that glitters is not gold - towards process discovery techniques with guarantees. In *Advanced Information Systems Engineering 2021*, volume 12751 of *LNCS*, pages 141–157. Springer, 2021.
- [39] M. L. van Eck, X. Lu, S. J. J. Leemans, and W. M. P. van der Aalst. PM2: A process mining project methodology. In *CAiSE 2015*, volume 9097 of *LNCS*, pages 297–313. Springer, 2015.

- [40] A. J. M. M. Weijters and J. T. S. Ribeiro. Flexible heuristics miner (FHM). In *CIDM 2011*, pages 310–317. IEEE, 2011.
- [41] B. R. van Wensveen. Estimation and analysis of the quality of event log samples for process discovery. Master’s thesis, Utrecht University, 2020. <https://dspace.library.uu.nl/handle/1874/400143>.